

Machine translation JP11338646

(19) **Publication country** Japan Patent Office (JP)
(12) **Kind of official gazette** Open patent official report (A)
(11) **Publication No.** JP,11-338646,A
(43) **Date of Publication** December 10, Heisei 11 (1999)
(54) **Title of the Invention** Disk array equipment
(51) **International Patent Classification (6th Edition)**
G06F 3/06 304
540

FI

G06F 3/06 304 B
540

Request for Examination Un-asking.

The number of claims 3

Mode of Application OL

Number of Pages 8

(21) **Application number** Japanese Patent Application No. 10-147552

(22) **Filing date** May 28, Heisei 10 (1998)

(71) **Applicant**

Identification Number 000005108

Name Hitachi, Ltd.

Address 4-6, Kanda Surugadai, Chiyoda-ku, Tokyo

(72) **Inventor(s)**

Name **** Katsumi

Address 2880, Kozu, Odawara-shi, Kanagawa-ken Inside of the Hitachi, Ltd. storage system operation division

(74) **Attorney**

Patent Attorney

Name Tsutsui Yamato

(57) **Abstract**

Technical problem Remote mirroring of disk array equipment is realized cheaply.

Means for Solution While connecting the disk drive 107 inside disk array equipment 100 to the drive interface 105 of a controller 101 through the fiber channel cable 106 and constituting FC_AL The external port 108 and the fiber channel cable 109 are minded. The disk drive 204 in the external disk unit 200 of a remote place It connects with the drive interface 105 of a controller 101, and FC_AL is constituted. The disk array of the disk drive 107 local by the common controller 101, Disaster recovery remote mirroring which stores the same data in multiplex between the disk arrays of the disk drive 204 of RIMOTO is realized.

(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号

特開平11-338646

(43)公開日 平成11年(1999)12月10日

(51)Int.Cl.⁶

G 0 6 F 3/06

識別記号

3 0 4

5 4 0

F I

G 0 6 F 3/06

3 0 4 B

5 4 0

審査請求 未請求 請求項の数3 O L (全 8 頁)

(21)出願番号 特願平10-147552

(22)出願日 平成10年(1998)5月28日

(71)出願人 000005108

株式会社日立製作所

東京都千代田区神田駿河台四丁目6番地

(72)発明者 饒崎 克巳

神奈川県小田原市国府津2880番地 株式会

社日立製作所ストレージシステム事業部内

(74)代理人 弁理士 筒井 大和

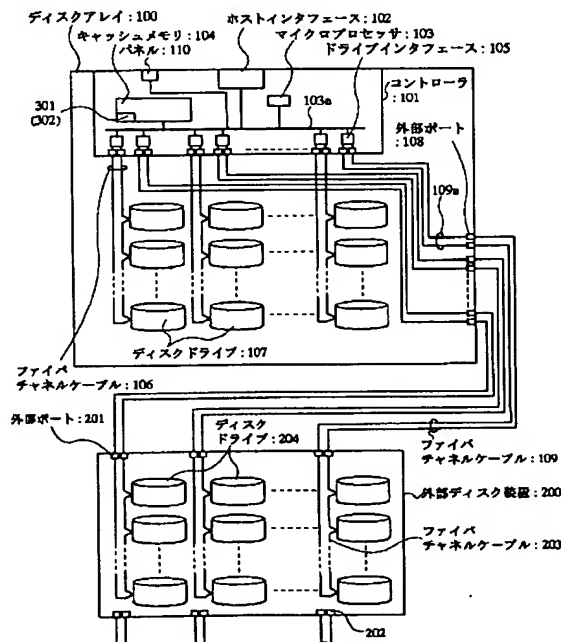
(54)【発明の名称】 ディスクアレイ装置

(57)【要約】

【課題】 ディスクアレイ装置の遠隔ミラーリングを安価に実現する。

【解決手段】 ディスクアレイ装置100の内部のディスクドライブ107を、ファイバチャネルケーブル106を介してコントローラ101のドライブインタフェース105に接続してFC_A_Lを構成するとともに、外部ポート108、ファイバチャネルケーブル109を介して、遠隔地の外部ディスク装置200内のディスクドライブ204を、コントローラ101のドライブインタフェース105に接続してFC_A_Lを構成し、共通のコントローラ101にて、ローカルのディスクドライブ107のディスクアレイと、リモートのディスクドライブ204のディスクアレイとの間で同一データを多重に格納するディザスタリカバリ遠隔ミラーリングを実現する。

図 1



【特許請求の範囲】

【請求項 1】 複数の単体ドライブと、前記単体ドライブの動作を制御するコントローラとを含むディスクアレイ装置であって、

前記単体ドライブは、前記ディスクアレイ装置の内部に配置され、前記コントローラと複数の第 1 のインタフェースにて接続される第 1 の単体ドライブと、前記ディスクアレイ装置の外部に配置され、前記ディスクアレイ装置に設けられた外部接続ポートを介して、前記コントローラと複数の第 2 のインタフェースにて接続される第 2 の単体ドライブと、

を含むことを特徴とするディスクアレイ装置。

【請求項 2】 複数の単体ドライブと、前記単体ドライブの動作を制御するコントローラとを含むディスクアレイ装置であって、

前記単体ドライブは、前記ディスクアレイ装置の内部に配置され、前記コントローラと複数のインタフェースにて接続される第 1 の単体ドライブと、

前記ディスクアレイ装置の外部に配置され、前記ディスクアレイ装置に設けられた外部接続ポートを介して、前記インタフェースに接続される第 2 の単体ドライブと、を含むことを特徴とするディスクアレイ装置。

【請求項 3】 請求項 1 または 2 記載のディスクアレイ装置において、

前記インタフェースは、ファイバチャネル・アービトラレーテッド・ループ（F C _ A L）であり、前記コントローラは、同一のデータを、前記第 1 および第 2 の単体ドライブに多重に格納することを特徴とするディスクアレイ装置。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、ディスクアレイ技術に関し、特に、ディスクアレイにおけるデータ多重化技術等に適用して有効な技術に関する。

【0002】

【従来の技術】コンピュータシステムの規模は、年々大規模化してきており、それに伴い扱うデータ量も増大してきている。このデータ量増大に伴い、大容量記憶システムとして耐故障性を考慮し、ディスクアレイ装置を使用する傾向に有る。これはディスクドライブが突然起こす故障に対して、データ保証を行う技術である。このディスクアレイ技術の中の一つに、ミラーリングという方法が有る。ミラーリングとは、複数台のディスクドライブのそれぞれに対し、同じデータを二重に書き込む手法である。このミラーリングの手段には次のようなものが存在する。

【0003】（1）ディスクアレイ装置を R A I D で使用する。

【0004】（2）ソフトウェア R A I D を使用し、別々のディスクに同じデータを書き込む。

【0005】（3）たとえば、特開平 8 - 2 7 2 6 6 6 号公報に開示された技術のように、ネットワークを使用し、ローカルのディスクアレイ装置に書き込んだデータをリモートのディスクアレイ装置にも書き込む。

【0006】前記（1）、（2）は、ローカルにてデータの二重化を行うものである。これは、単体のディスクドライブの故障によるデータ喪失防止を行うことは可能であるが、地震や火事等といった災害が発生した場合、二重化したディスクドライブが、共に破壊される可能性があり、データを二重化していてもデータを喪失する可能性が非常に高い。

【0007】一方、（3）はネットワークを使用しているため、遠隔地に二重化したデータが存在するので、災害によりローカルのディスクドライブが全て破壊されたとしても、遠隔地のディスクからデータ復旧が可能であるが、遠隔地に対してもディスクアレイ装置を必要とするため、コストがかかる。

【0008】

【発明が解決しようとする課題】上述のように、ディスクアレイ装置のミラーリングは、アレイ内のディスクにて行うため、災害発生時にミラーリングしたデータも同時に破壊され、データ喪失に至る、という技術的課題がある。

【0009】また、遠隔ミラーリングをディスクアレイ装置にて行う場合、中央処理装置、もしくは、ディスクアレイ装置がネットワークを使用して、遠方のディスク装置にデータを転送するが、その際は遠方のディスク装置に対しても、別のディスクアレイ装置が必要となり、ディスクアレイにおけるデータ二重化のコストが高くなる、という技術的課題もある。

【0010】本発明の目的は、災害等によるデータ喪失を確実に防止することが可能なディスクアレイ技術を提供することにある。

【0011】本発明の他の目的は、安価に、遠隔ミラーリングを実現することが可能なディスクアレイ技術を提供することにある。

【0012】

【課題を解決するための手段】本発明では、ディスクアレイ装置において、装置内部に配置されたローカルの第 1 の単体ドライブ、およびディスクアレイ装置の外部に設けられたリモートの第 2 の単体ドライブを、アレイコントローラに対してファイバチャネル（F C _ A L : F i b r e C h a n n e l - A r b i t r a t e d L o o p）等のインタフェースにて接続するものである。

【0013】より具体的には、一例として、ディスクアレイ装置にファイバチャネル用制御ポートを設け、遠隔地にあるディスク装置（ドライブ群）にも、ファイバチャネル用制御ポートを設けて、お互いをファイバチャネルプロトコルを使用して接続する。

【0014】ファイバチャネル用制御ポート数、および

そのポートに FC__AL 接続するドライブの数は、ローカルのディスクアレイのポート数、およびポートに付属のドライブ数と同数とし、ローカルのディスクアレイの構成と同様にする。

【0015】ファイバチャネルプロトコルにより、遠隔地のドライブに、ローカルのディスクアレイ装置を接続可能であり、かつアレイコントローラ等の特別な制御装置を遠隔地のディスク装置に持つ必要がなく、安価にミラーリングを行うことが可能である。

【0016】ファイバチャネルプロトコルによりお互いを接続することで、ネットワークを使用せず遠方にあるミラーリング用ドライブに接続可能であり、かつ特別な制御装置をミラーリング用ドライブ群側に持つ必要がなくなり、安価に遠隔ミラーリングを行うことが可能となる。

【0017】

【発明の実施の形態】以下、本発明の実施の形態を図面を参照しながら詳細に説明する。

【0018】ミラーリングとは、データの二重化のことである。本実施の形態のディスクアレイ装置では、主ドライブと副ドライブを定義する。上位装置から発行されたデータは、主ドライブに書き込むと同時に副ドライブにも書き込まれる。ドライブからのデータ読み出しは、基本的に主ドライブから行うが、主ドライブが何らかの故障によりアクセス不能となった場合、副ドライブからデータを読み出す。このような方法で、ディスクドライブの障害によるデータロス防止を行うことをミラーリングといい、特に、ディスクアレイでは RAID1 と呼ばれる。この RAID1 は、ディスクアレイ内部のディスク群を主および副として使用するため、火事や地震といった災害などによりディスクアレイ装置そのものが壊れた場合データを失う。

【0019】そこで、本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムは、この技術的課題を解決するために、以下のように、ファイバチャネルプロトコルを使用し、遠隔地に設置されたディスクドライブに対してミラーリングを行うものである。

【0020】図1は、本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムの構成の一例を示す概念図である。

【0021】本実施の形態のディザスタリカバリ遠隔ミラーリングシステムは、ディスクアレイ装置100と外部ディスク装置200から構成される。

【0022】ディスクアレイ装置100は、コントローラ101、外部ポート108、ディスクドライブ107、ファイバチャネルケーブル106から構成される。

【0023】コントローラ101は、ディザスタリカバリ遠隔ミラーリングシステム全体を制御し、上位装置であるホストとディスクアレイ装置100の間のプロトコ

ル制御を司るホストインタフェース102、装置全体の制御を行うマイクロプロセッサ103、ホストからのデータ、および後述の図2に例示されるような制御情報を保持するキャッシュメモリ104、ディスクドライブ107、204とのプロトコル制御を行うドライブインタフェース105から構成される。マイクロプロセッサ103に対して、それ以外の各部はシステムバス103aを介して接続され、マイクロプロセッサ103の制御の下で動作する。

【0024】ディスクドライブ107は、ファイバチャネルケーブル106を使用し、コントローラ101のドライブインタフェース105と接続する。

【0025】外部ポート108は、外部ディスク装置200とファイバチャネルケーブル109を使用して接続するためのインタフェースであり、ファイバチャネルケーブル109aを介して、コントローラ101のドライブインタフェース105と接続されている。

【0026】外部ディスク装置200は、外部ポート201、外部ポート202、複数のディスクドライブ204、ファイバチャネルケーブル203から構成される。外部ポート201は、ディスクアレイ装置100とディスクドライブ204間のインタフェースであり、ディスクドライブ204はファイバチャネルケーブル203を使用してディスクアレイ装置100と接続される。

【0027】なお、外部ポート202は、増設用に用いられ、たとえばデータを3重以上に多重化する場合に、たとえば外部ディスク装置200と同様の構成を持つ外部ディスク装置が、必要に応じて接続される。

【0028】次に、本実施の形態にて用いられる制御情報の一例について示す。図2は本実施の形態で使用する記憶管理体系（制御情報）の一例を示す概念図である。本実施の形態では、一例として図2に例示されるドライブ管理テーブル301と、データ管理テーブル302とを用いる。ドライブ管理テーブル301は、主ドライブおよび副ドライブの各々について、ドライブの数分のエントリを持つ。

【0029】ミラーリングを行う主ドライブ（ディスクドライブ107）と副ドライブ（ディスクドライブ204）は、それぞれ別々のファイバチャネルケーブル106、およびファイバチャネルケーブル109で、コントローラ101のドライブインタフェース105に接続され、独立な FC__AL ループを構成する。

【0030】ミラーリングを制御するに際し、まず主ドライブと副ドライブの対応を取る必要が有る。ファイバチャネルにおけるデバイスは、それぞれ世界中で一意的 WWN (World Wide Name) を持つ。ディスクドライブ107、ディスクドライブ204の WWN は、 FC__AL ループ初期化時に、コントローラ101が認識可能であるが、この時、あるディスクドライブが、ディスクアレイ装置100に属しているのか、外部ディスク装置2

00に属しているのか、認識不可能である。また外部ディスク装置200は外部ポート202を使用した場合に、複数台接続可能なので、どの外部ディスク装置200にあるのかも不明である。そこで、ディスクアレイ装置100のパネル110からディスクアレイ装置100に属しているディスクのWWNを設定し、ドライブ管理テーブル301のWWN301-1-cおよびWWN301-2-cに登録する。

【0031】また、事前にパネル110を使用し、ループID301-1-b、ループID301-2-bをそれぞれ設定することによって、主側に使用するループと副側に使用するループの登録を行う。

【0032】D_ID301-1-a、D_ID301-2-aはPORT_IDでファイバチャネルでデバイスを特定するのに使用される。これはFC_AL初期化時に動的に決定される。

【0033】ドライブ正常フラグ301-1-e、ドライブ正常フラグ301-2-eは、そのドライブが正常か異常かを示す。異常の場合はそのドライブに対し、I/Oを行わない。副ドライブへのポインタ301-1-d、主ドライブへのポインタ301-2-dはミラーリングを行っているディスクの互いの対応ドライブを認識するためのエントリへのポインタである。

【0034】上位装置からの書込データは、たとえば図2に例示されるデータ管理テーブル302を使用して管理する。データへのポインタ302-aはデータの存在するキャッシュのアドレスである。D_ID302-bはデータの書き込む位置にあるドライブのD_IDである。LBA302-cはデータを反映させるディスク内のアドレスである。主書込未反映フラグ302-dは、このデータが主側のディスクドライブ107に対し既に反映されたかどうかを示す。同様に、副書込未反映フラグ302-eは、このデータが副側のディスクドライブ204に対し反映されたか否かを示す。

【0035】以下、本実施の形態の作用の一例について説明する。

【0036】図3のフローチャートにて、本実施の形態の遠隔ミラーリングの処理の一例を遠隔ミラーリングの制御法の一例として示す。

【0037】上位装置から転送されたコマンドはディスクアレイ制御を行うコントローラ101により解析され、WRITE/READ等が判別される(ステップ401)。

【0038】WRITEコマンドの場合、上位装置から受領したコマンドに対するデータが転送され、このデータをキャッシュメモリ104に保持する(ステップ402)。

【0039】このとき、WRITEコマンドのロジカルブロックアドレス(LBA)およびレングス(データ長)からローカルのディスクアレイ装置100上のディ

スクドライブ107のD_IDとLBA、レングスを算出し、データ管理テーブル302を設定する(ステップ403)。

【0040】次に、キャッシュメモリ104上のデータをディスクドライブ107(204)に書き込む。キャッシュメモリ104上のデータをディスクドライブ107(204)に書き込むタイミングは、たとえば、タイマにより時間を監視し、一定の時間がきたらディスクドライブ107(204)に反映させる、もしくはキャッシュメモリ104の使用エリアの容量に閾値を使用し、閾値を超えたらディスクドライブ107(204)に反映させる、もしくは、I/Oの負荷が大きくないときにディスクドライブ107(204)に反映させるといった形で非同期に行い、不要なディスクドライブ107(204)へのアクセスを減らす。

【0041】書込契機になると、まず、データ管理テーブル302のD_ID302-bをタグとしてドライブ管理テーブル301のD_ID301-1-a、副ドライブへのポインタ301-1-dから、副ドライブの有無を確認し(ステップ404)、ない場合には、主側のディスクドライブ107のみにデータ書込を実行し(ステップ405)、キャッシュメモリ104の領域を解放した後(ステップ406)、ホストにWRITE完了を応答する(ステップ407)。

【0042】一方、副ディスクドライブが存在する場合には、ローカルのディスクドライブ107および対応するミラーリングディスクに対し同時に書込処理行う。

【0043】すなわち、データ管理テーブル302のD_ID302-bをタグとしてドライブ管理テーブル301のD_ID301-1-a、副ドライブへのポインタ301-1-dから、対応する主/副ドライブを特定した後(ステップ408)、データ管理テーブル302の主書込未反映フラグ302-d、副書込未反映フラグ302-eを共にONにし(ステップ409)、その後、主/副ドライブの両方にデータ書込命令を発行する(ステップ410)。

【0044】主/副ドライブの各々では、それぞれWRITE(データ反映)完了を監視し(ステップ411、ステップ413)、ディスクドライブ107、ディスクドライブ204の各々ではデータを反映すると、反映した側のディスクドライブに対するデータ未反映フラグをOFFとする(ステップ412、ステップ414)。

【0045】その際、対応する他方のディスクのデータ未反映フラグを確認し、OFFとなっていればキャッシュ領域として再びデータを書き込むことを可能とし、ONのままであれば、これがOFFになるまでこのキャッシュ領域の使用を不可とし、ミラーリングが未反映となることを防止する(ステップ415)。

【0046】主/副のディスクドライブ107および204の両方にデータの反映が完了すると、キャッシュメ

モリ 104 の領域を解放した後（ステップ 406）、ホストに WRITE 完了を応答する（ステップ 407）。

【0047】一方、READ コマンドの場合、基本的に主側のディスクドライブ 107 からデータを取得するが、この主側のディスクドライブ 107 が故障していた場合、データを取得できないため、この時は対応する副側のディスクドライブ 204 からデータを取得する。

【0048】すなわち、図 4 のフローチャートに例示されるように、まず、コマンド受領時にディスクアレイ装置上のディスクの D_ID と LBA、レングスを算出した後（ステップ 501）、アクセス対象の主側のディスクドライブ 107 が正常か否かを、ドライブ管理テーブル 301 の当該ドライブのドライブ正常フラグ 301-1-e を参照して判別し（ステップ 502）、正常な場合には、主側のディスクドライブ 107 からデータを読み出して（ステップ 503）、ホスト側に転送する（ステップ 504）。

【0049】異常の場合には、データ管理テーブル 302 の D_ID 302-b をタグとしてドライブ管理テーブル 301 の D_ID 301-1-a、副ドライブへのポインタ 301-1-d から、対応する主/副ドライブを特定したのち（ステップ 505）、特定された副側のディスクドライブ 204 からデータを読み出して（ステップ 506）、ホスト側に転送する（ステップ 504）。

【0050】このように、本実施の形態のディザスタリカバリ遠隔ミラーリングによれば、ディスクアレイ装置 100 内のローカルな複数のディスクドライブ 107 と、遠隔地の外部ディスク装置 200 に設けられた複数のディスクドライブ 204 の各々を、それぞれ別々のファイバチャネルケーブル 106、およびファイバチャネルケーブル 109 で、コントローラ 101 のドライブインタフェース 105 に接続して、独立な FC_AL ループを構成し、対応するディスクドライブ 107 とディスクドライブ 204 に同じデータを格納する遠隔ミラーリングを行うので、中央処理装置等のホストからディスクアレイ装置 100 へのライトデータのコピーを遠隔地の外部ディスク装置 200 に保持することが可能となり、災害発生時にディスクアレイ装置 100 が破壊されても遠隔地の外部ディスク装置 200 にてデータ復旧が可能となる。

【0051】また、ファイバチャネルプロトコルを使用して、ディスクアレイ装置 100 内のコントローラ 101 と、遠隔地の外部ディスク装置 200 と接続しているため、遠隔地側にある外部ディスク装置 200 はアレイコントローラを備える必要がなく、安価に遠隔ミラーリングを行うことが可能となる。

【0052】なお、図 1 に例示した構成では、ディスクアレイ装置 100 内のローカルのディスクドライブ 107 と、外部ディスク装置 200 のディスクドライブ 204 とは、ファイバチャネルケーブル 106 の系列の FC

__AL と、ファイバチャネルケーブル 109 a、外部ポート 108、ファイバチャネルケーブル 109 からなる系列の FC__AL とに独立に接続されていたが、これに限らず、たとえば、図 5 に例示されるように、ディスクドライブ 107 が接続されるファイバチャネルケーブル 106 a を、外部ポート 108 に接続し、この外部ポート 108 を介して、ファイバチャネルケーブル 109 と接続することにより、主側のディスクドライブ 107 と、外部ディスク装置 200 内の副側のディスクドライブ 204 とが、同一系列の FC__AL に属する構成としてもよい。

【0053】そして、ローカルのディスクドライブ 107 と、対応する外部のディスクドライブ 204 に同じデータを二重に格納する遠隔ミラーリングを行う。

【0054】この図 5 の変形例の場合には、図 2 に例示されたドライブ管理テーブル 301 で、主/副のディスクドライブの各々が帰属する FC__AL を区別するための、ループ ID 301-1-b、ループ ID 301-2-b のデータは不要であり、制御情報が簡素化されるとともに、ディスクアレイ装置 100 内のコントローラ 101 におけるドライブインタフェース 105 の数を減らすことができ、ハードウェア構成が簡素化される、という利点がある。

【0055】以上本発明者によってなされた発明を実施の形態に基づき具体的に説明したが、本発明は前記実施の形態に限定されるものではなく、その要旨を逸脱しない範囲で種々変更可能であることはいうまでもない。

【0056】たとえばインタフェースとしては、ファイバチャネルに限らず、遠距離のデータ授受が可能なインタフェースであれば、他のインタフェースでもよい。

【0057】

【発明の効果】本発明のディスクアレイ装置によれば、災害等によるデータ喪失を確実に防止することができ、という効果が得られる。

【0058】また、安価に、遠隔ミラーリングを実現することができる、という効果が得られる。

【図面の簡単な説明】

【図 1】本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムの構成の一例を示す概念図である。

【図 2】本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムで使用する制御情報の一例を示す概念図である。

【図 3】本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムの作用の一例を示すフローチャートである。

【図 4】本発明のディスクアレイ装置の一実施の形態であるディザスタリカバリ遠隔ミラーリングシステムの作用の一例を示すフローチャートである。

【図 5】本発明のディスクアレイ装置の一実施の形態で

9

あるディザスタリカバリ遠隔ミラーリングシステムの変形例を示す概念図である。

【符号の説明】

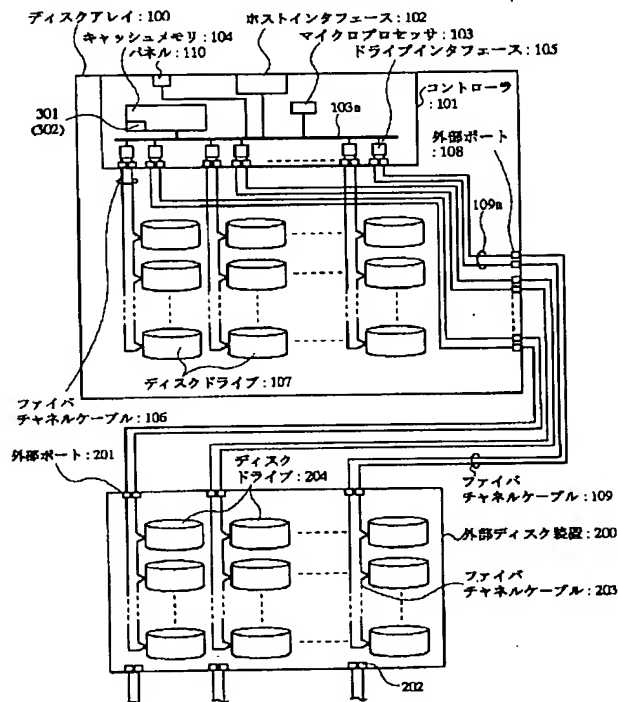
100…ディスクアレイ装置、101…コントローラ、102…ホストインタフェース、103…マイクロプロセッサ、104…キャッシュメモリ、105…ドライブインタフェース、106…ファイバチャネルケーブル（第1のインタフェース）、106a…ファイバチャネルケーブル（インタフェース）、107…ディスクドライブ（第1の単体ドライブ）、108…外部ポート、109…ファイバチャネルケーブル（第2のインタフェース）、109a…ファイバチャネルケーブル、110…パネル、200…外部ディスク装置、201…外部ポ

10

ポート、202…外部ポート、203…ファイバチャネルケーブル、204…ディスクドライブ（第2の単体ドライブ）、301…ドライブ管理テーブル、301-1-a…D_ID、301-1-b…ループID、301-1-c…WWN、301-1-d…副ドライブへのポインタ、301-1-e…ドライブ正常フラグ、301-2-a…D_ID、301-2-b…ループID、301-2-c…WWN、301-2-d…主ドライブへのポインタ、301-2-e…ドライブ正常フラグ、302…データ管理テーブル、302-a…キャッシュメモリ上のデータへのポインタ、302-b…D_ID、302-c…LBA、302-d…主書込未反映フラグ、302-e…副書込未反映フラグ。

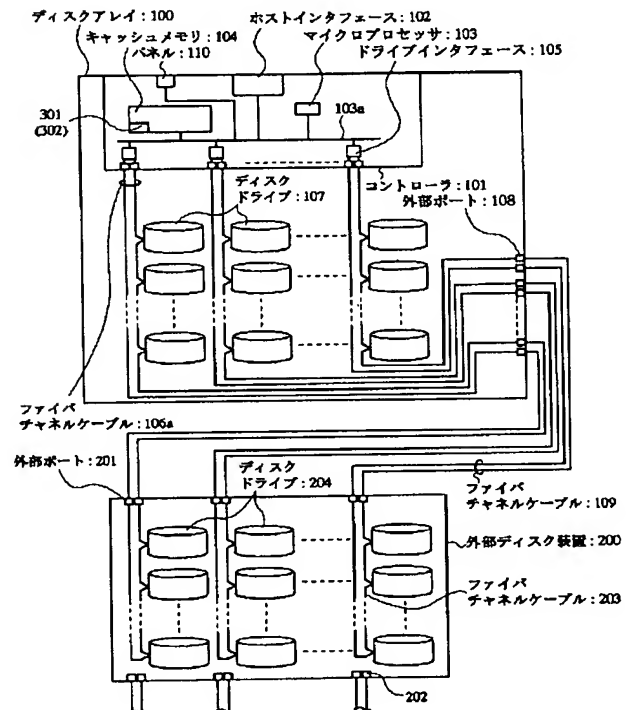
【図1】

図 1



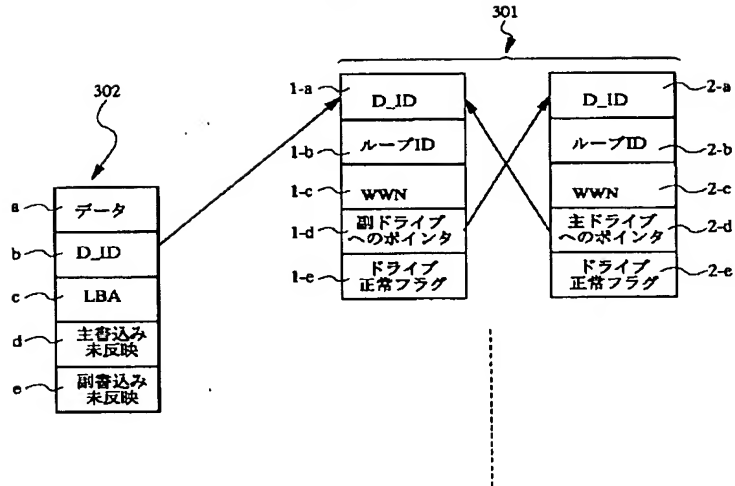
【図5】

図 5



【図 2】

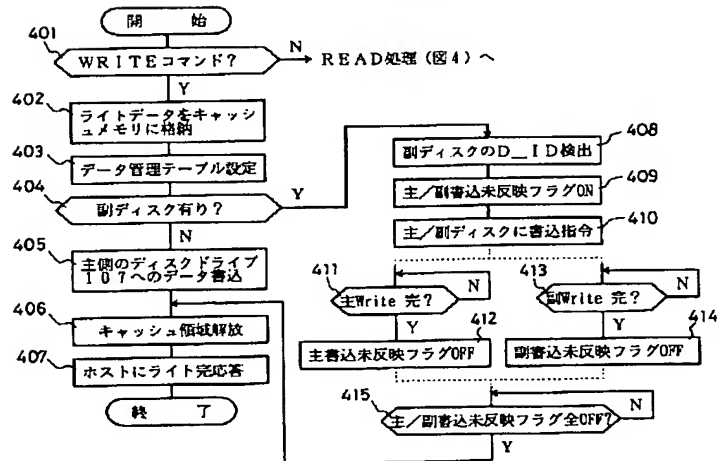
図 2



【図 3】

図 3

WRITEコマンド処理



【図 4】

図 4

